



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

Mean field for Markov Decision Processes: from Discrete to Continuous Optimization

Nicolas Gast — Bruno Gaujal — Jean-Yves Le Boudec

N° 7239 — version 3

initial version Avril 2010 — revised version May 2011

Thème NUM

A large blue rectangle occupies the lower half of the page. Overlaid on it is a large, light gray stylized 'R' logo. To the right of the 'R', the words 'Rapport de recherche' are written in a white serif font, with 'Rapport' on the top line and 'de recherche' on the bottom line. A horizontal gray brushstroke underline is positioned below the text.

*Rapport
de recherche*

Mean field for Markov Decision Processes: from Discrete to Continuous Optimization

Nicolas Gast, Bruno Gaujal, Jean-Yves Le Boudec

Thème NUM — Systèmes numériques
Équipe-Projet MESCAL

Rapport de recherche n° 7239 — version 3 — initial version Avril 2010 — revised version May 2011 — 25 pages

Abstract: We study the convergence of Markov decision processes, composed of a large number of objects, to optimization problems on ordinary differential equations. We show that the optimal reward of such a Markov decision process, which satisfies a Bellman equation, converges to the solution of a continuous Hamilton-Jacobi-Bellman (HJB) equation based on the mean field approximation of the Markov decision process. We give bounds on the difference of the rewards and an algorithm for deriving an approximating solution to the Markov decision process from a solution of the HJB equations. We illustrate the method on three examples pertaining, respectively, to investment strategies, population dynamics control and scheduling in queues. They are used to illustrate and justify the construction of the controlled ODE and to show the advantage of solving a continuous HJB equation rather than a large discrete Bellman equation.

Key-words: Mean Field, Hamilton-Jacobi-Bellman, Optimal Control, Markov Decision Process

Modèles Champ Moyen et Processus de Décision Markovien: de l'optimisation discrète à l'optimisation continue.

Résumé : Ce document étudie la convergence de processus de décision markoviens composés d'un grand nombre d'objets vers des problèmes d'optimisation sur des équations différentielles. Nous montrons que le gain optimal du processus de décision converge vers la solution d'une équation continue de type "Hamilton-Jacobi-Bellman". La preuve utilise à la fois des outils classiques des modèles champs moyens et différents nouveaux couplages entre les modèles discrets et continus qui permettent de donner des bornes explicites. La méthode est ensuite illustrée par trois exemples concernant des stratégies d'investissement, du contrôle de dynamiques de population et un problème d'allocation de ressources.

Mots-clés : Champ Moyen, Hamilton-Jacobi-Bellman, Contrôle Optimal, Processus de Décision Markovien

1 Introduction

In this paper we study dynamic optimization problems on Markov decision processes composed of a large number of interacting objects.

Consider a system of N objects evolving in a common environment. At each time step, objects change their state randomly according to some probability kernel Γ^N . This kernel depends on the number of objects in each state, as well as on the decisions of a centralized controller. Our goal is to study the behavior of the controlled system when N becomes large.

Several papers investigate the asymptotic behavior of such systems, but without controllers. For example, in [2, 19], the authors show that under mild conditions, as N grows, the system converges to a deterministic limit. The limiting system can be of two types, depending on the intensity $I(N)$ (the intensity is the probability than an object changes its state between two time steps). If $I(N) = O_{N \rightarrow \infty}(1)$, the system converges to a dynamical system in discrete time [19]. If $I(N)$ goes to 0 as N grows, the limiting system is a continuous time dynamical system and can be described by ordinary differential equations (ODEs).

Contributions

Here, we consider a Markov decision process where at each time step, a central controller chooses an action from a predefined set that will modify the dynamics of the system the controller receives a reward depending on the current state of the system and on the action. The goal of the controller is to maximize the expected reward over a finite time horizon. We show that when N becomes large this problem converges to an optimization problem on an ordinary differential equation.

More precisely, we focus on the case where the Markov decision process is such that its empirical occupancy measure is also Markov; this occurs when the system consists of many interacting objects, the objects can be observed only through their state and the system evolution depends only on the collection of all states. We show that the optimal reward converges to the optimal reward of the mean field approximation of the system, which is given by the solution of an HJB equation. Furthermore, the optimal policy of the mean field approximation is also asymptotically optimal in N , for the original discrete system. Our method relies on bounding techniques used in stochastic approximation and learning [4, 1]. We also introduce an original coupling method, where, to each sample path of the Markov decision process, we associate a random trajectory that is obtained as a solution of the ODE, i.e. the mean field limit, controlled by random actions.

This convergence result has an algorithmic by-product. Roughly speaking, when confronted with a large Markov decision problem, we can first solve the HJB equation for the associated mean field limit and then build a decision policy for the initial system that is asymptotically optimal in N .

Our results have two main implications. The first is to justify the construction of controlled ODEs as good approximations of large discrete controlled systems. This construction is given done without rigorous proofs. In Section 4.3.2 we illustrate this point with an example of malware infection in computer systems.

The second implication concerns the effective computation of an optimal control policy. In the discrete case, this is usually done by using dynamic programming for the finite horizon case or by computing a fixed point of the Bellman equation in the discounted case. Both approaches suffer from the curse of dimensionality, which makes them impractical when the state space is too large. In our context, the size of the state space is exponential in N , making the problem even more acute. In practice, modern supercomputers only allow us to tackle such optimal control problems when N is no larger than a few tens [20].

The mean field approach offers an alternative to brute force computations. By letting N go to infinity, the discrete problem is replaced by a limit Hamilton-Jacobi-Bellman equation that is deterministic where the dimensionality of the original system has been hidden in the occupancy measure. Solving the HJB equation numerically is sometimes rather easy, as in the examples in Sections 4.3.1 and 4.3.2. It provides a deterministic optimal policy whose reward with a finite (but large) number of objects is remarkably close to the optimal reward.

Related Work

Several papers in the literature are concerned with the problem of mixing the limiting behavior of a large number of objects with optimization.

In [6], the value function of the Markov decision process is approximated by a linearly parametrized class of functions and a fluid approximation of the MDP is used. It is shown that a solution of the HJB equation is a value function for a modification of the original MDP problem. In [25, 8], the curse of dimensionality of dynamic programming is circumvented by approximating the value function by linear regression. Here, we use instead a mean field limit approximation and prove asymptotic optimality in N of limit policy.

In [9], the authors also consider Markov decision processes with a growing number of objects, but when the intensity is $O(1)$. In their case, the optimization problem of the system of size N converges to a deterministic optimization problem in discrete time. In this paper however, we focus on the $o(1)$ case, which is substantially different from the discrete time case because the limiting system does not evolve in discrete time anymore.

Actually, most of the papers dealing with mean field limits of optimization problems over large systems are set in a game theory framework, leading to the concept of *mean field games* introduced in [18]. The objects composing the system are seen as N players of a game with distributed information, cost and control; their actions lead to a Nash equilibrium. To the best of our knowledge, the classic case with global information and centralized control has not yet been considered. Our work focuses precisely on classic Markov decision problems, where a central controller (our objects are passive), aims at minimizing a global cost function.

For example, a series of papers by M. Huang, P.E. Caines and P. Malhamé such as [11, 12, 13, 14] investigate the behavior of systems made of a large number of objects under *distributed* control. They mostly investigate Linear-Quadratic- Gaussian (LQG) dynamics and use the fact that, here, the solution can be given in closed form as a Riccati equation to show that the limit satisfies a Nash fixed point equation. Their more general approach uses the Nash Equivalence Certainty principle introduced in [11]. The limit equilibrium could or could not be a global optimal. Here, we consider the general case where the dynamics and the cost may be arbitrary (we do not assume LQG Dynamics) so that the optimal policy is not given in closed form. The main difference with their approach comes from the fact that we focus instead on centralized control to achieve a global optimum. The techniques to prove convergence are rather different. Our proofs are more in line with classic mean field arguments and use stochastic approximation techniques.

Another example is the work of Tembiné and others [23, 24], on the limits of games with many players. The authors provide conditions under which the limit when the number of players grows to infinity commutes with the fixed point equation satisfied by a Nash equilibrium. Again, our investigation solves a different problem and focuses on the centralized case. In addition, our approach is more algorithmic; we construct two intermediate systems: one with a finite number of objects controlled by a limit policy and one with a limit system controlled by a stochastic policy induced by the finite system.

Structure of the paper

The rest of the paper is structured as follows. In Section 2 we give definitions, some notation and hypotheses. In Section 3 we describe our main theoretical. In Section 4 we describe our resulting algorithm and illustrate the application of our method with a few examples. The details of all proofs are in Section 5 and Section 6 concludes the paper.

2 Notations and Definitions

2.1 System with N Objects

We consider a system composed of N *objects*. Each object has a state from the finite set $SS = \{1 \dots S\}$. Time is discrete and the state of the object n at step $k \in \mathbb{N}$ is denoted $X_n^N(k)$. The state

of the system at time k is $X^N(k) \stackrel{\text{def}}{=} (X_1^N(k) \dots X_N^N(k))$. For all $i \in SS$, we denote by $M^N(k)$ the empirical measure of the objects $(X_1^N(k) \dots X_N^N(k))$ at time k :

$$M^N(k) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \delta_{X_n^N(k)}, \quad (1)$$

where δ_x denotes the Dirac measure in x . $M^N(k)$ is a probability measure on SS and its i th component $M^N(k)[i]$ denotes the proportions of objects in state i at time k (also called the occupancy measure): $M^N(k)[i] = \frac{1}{N} \sum_{n=1}^N \mathbf{1}_{X_n^N(k)=i}$.

The system $(X^N(k))_{k \in \mathbb{N}}$ is a Markov process once the sequence of the actions taken by the controller is fixed. Let Γ^N be the transition kernel, namely Γ^N is a mapping $SS^N \times SS^N \times \mathcal{A} \rightarrow [0, 1]$, where \mathcal{A} is the set of possible actions, such that for every $x \in SS^N$ and $a \in \mathcal{A}$, $\Gamma^N(x, \cdot, a)$ is a probability distribution on SS^N and further, if the controller takes the action $A^N(k)$ at time t and the system is in state $X^N(k)$, then:

$$\mathcal{P}(X^N(k+1) = y_1 \dots y_N | X^N(k) = x_1 \dots x_N, A^N(k) = a) = \Gamma^N(x_1 \dots x_N, y_1 \dots y_N, a) \quad (2)$$

We assume that

(A0) Objects are observable only through their states

in particular, the controller can observe the collection of all states X_1^N, X_2^N, \dots , but not the identities $n = 1, 2, \dots$. This assumption is required for mean field convergence to occur. In practice, it means that we need to put into the object state any information that is relevant to the description of the system.

Assumption (A0) translates into the requirement that the kernel be invariant by object relabeling. Formally, let \mathfrak{S}^N be the set of permutations of $\{1, 2, \dots, N\}$. By a slight abuse of notation, for $\sigma \in \mathfrak{S}^N$ and $x \in SS^N$ we also denote with $\sigma(x)$ the collection of object states after the permutation, i.e. $\sigma(x) \stackrel{\text{def}}{=} (x_{\sigma^{-1}(1)} \dots x_{\sigma^{-1}(N)})$. The requirement is that

$$\Gamma^N(\sigma(x), \sigma(y), a) = \Gamma^N(x, y, a) \quad (3)$$

for all $x, y \in SS^N$, $\sigma \in \mathfrak{S}^N$ and $a \in \mathcal{A}$. A direct consequence, shown in Section 5, is:

Theorem 1. *For any given sequence of actions, the process $M^N(t)$ is a Markov chain*

2.2 Action, Reward and Policy

At every time k , a centralized controller chooses an action $A^N(k) \in \mathcal{A}$ where \mathcal{A} is called the action set. (\mathcal{A}, d) is a compact metric space for some distance d . The purpose of Markov decision control is to compute optimal *policies*. A policy $\pi = (\pi_0, \pi_1, \dots, \pi_k, \dots)$ is a sequence of decision rules that specify the action at every time instant. The policy π_k might depend on the sequence of past and present states of the process X^N , however, it is known that when the state space is finite, the action set compact and the kernel and the reward are continuous, there exists a deterministic Markovian policy which is optimal (see Theorem 4.4.3 in [21]). This implies that we can limit ourselves to policies that depend only on the current state $X^N(k)$.

Further, we assume that the controller can only observe object states. Therefore she cannot make a difference between states that result from object relabeling, i.e. the policy depends on $X^N(k)$ in a way that is invariant by permutation. By Lemma 2 in Section 5.2, it depends on $M^N(k)$ only. Thus, we may assume that, for every k , π_k is a function $\mathcal{P}(SS) \rightarrow \mathcal{A}$. Let $M_\pi^N(k)$ denotes the occupancy measure of the system at time k when the controller applies policy π .

If the system has occupancy measure $M^N(k)$ at time k and if the controller chooses the action $A^N(k)$, she gets an *instantaneous reward* $r^N(M^N(k), A^N(k))$. The expected value over a finite-time horizon $[0; H^N]$ starting from m_0 when applying the policy π is defined by

$$V_\pi^N(m) \stackrel{\text{def}}{=} \mathbb{E} \left(\sum_{k=0}^{\lfloor H^N \rfloor} r^N(M_\pi^N(k), \pi(M_\pi^N(k))) \middle| M_\pi^N(0) = m \right) \quad (4)$$

The goal of the controller is to find an optimal policy that maximizes the expected value. We denote by $V_*^N(m)$ the optimal value when starting from m :

$$V_*^N(m) = \sup_{\pi} V_{\pi}^N(m) \quad (5)$$

2.3 Scaling Assumptions

If at some time k , the system has occupancy measure $M^N(k) = m$ and the controller chooses action $A^N(k) = a$, the system goes into state $M^N(k+1)$ with probabilities given by the kernel $Q^N(M^N(k), A^N(k))$. The expectation of the difference between $M^N(k+1)$ and $M^N(k)$ is called the *drift* and is denoted by $F^N(m, a)$:

$$F^N(m, a) \stackrel{\text{def}}{=} \mathbb{E} [M^N(k+1) - M^N(k) | M^N(k) = m, A^N(k) = a]. \quad (6)$$

In order to study the limit with N , we assume that F^N goes to 0 at speed $I(N)$ when N goes to infinity and that $F^N/I(N)$ converges to a Lipschitz continuous function f . More precisely, we assume that there exists a sequence $I(N) \in (0; 1)$, $N = 1, 2, 3, \dots$, called the *intensity* of the model with $\lim_{N \rightarrow \infty} I(N) = 0$ and a sequence $I_0(N)$, $N = 1, 2, 3, \dots$, also with $\lim_{N \rightarrow \infty} I_0(N) = 0$ such that for all $m \in \mathcal{P}(SS)$ and $a \in \mathcal{A}$: $\left| \frac{1}{I(N)} F^N(m, a) - f(m, a) \right| \leq I_0(N)$. In a sense, $I(N)$ represents the order of magnitude of the number of objects that change their state within one unit of time.

The change of $M^N(k)$ during a time step is of order $I(N)$. This suggests a rescaling of time by $I(N)$ to obtain an asymptotic result. We define the continuous time process $(\hat{M}^N(t))_{t \in \mathbb{R}^+}$ as the affine interpolation of $M^N(k)$, rescaled by the intensity function, i.e. \hat{M}^N is affine on the intervals $[kI(N), (k+1)I(N)]$, $k \in \mathbb{N}$ and

$$\hat{M}^N(kI(N)) = M^N(k).$$

Similarly, \hat{M}_{π}^N denotes the affine interpolation of the occupancy measure under policy π . Thus, $I(N)$ can also be interpreted as the duration of the time slot for the system with N objects.

We assume that the time horizon and the reward per time slot scale accordingly, i.e. we impose

$$\begin{aligned} H^N &= \left\lfloor \frac{T}{I(N)} \right\rfloor \\ r^N(m, a) &= I(N)r(m, a) \end{aligned}$$

for every $m \in \mathcal{P}(SS)$ and $a \in \mathcal{A}$ (where $\lfloor x \rfloor$ denotes the largest integer $\leq x$).

2.4 Limiting System (Mean Field Limit)

We will see in Section 3 that as N grows, the stochastic system \hat{M}_{π}^N converges to a deterministic limit m_{π} , the mean field limit. For more clarity, all the stochastic variables (*i.e.*, when N is finite) are in uppercase and their limiting deterministic values are in lowercase.

An action function $\alpha : [0; T] \rightarrow \mathcal{A}$ is a piecewise Lipschitz continuous function that associates to each time t an action $\alpha(t)$. Note that action functions and policies are different in the sense that action functions do not take into account the state to determine the next action. For an action function α and an initial condition m_0 , we consider the following ordinary integral equation for $m(t)$, $t \in \mathbb{R}^+$:

$$m(t) - m(0) = \int_0^t f(m(s), \alpha(s)) ds. \quad (7)$$

(This equation is equivalent to an ODE, but is easier to manipulate in integral form. In the rest of the paper, we make a slight abuse of language and refer to it as an ODE). Under the foregoing assumptions on f and α , this equation satisfies the Cauchy Lipschitz condition and therefore has a

unique solution once the initial condition $m(0) = m_0$ is fixed. We call ϕ_t , $t \in \mathbb{R}^+$, the corresponding semi-flow, i.e.

$$m(t) = \phi_t(m_0, \alpha) \quad (8)$$

is the unique solution of Eq.(7).

As for the system with N objects, we define $v_\alpha(m_0)$ as the value of the limiting system over a finite horizon $[0; T]$ when applying the action function α and starting from $m(0) = m_0$:

$$v_\alpha(m_0) \stackrel{\text{def}}{=} \int_0^T r(\phi_s(m_0, \alpha), \alpha(s)) ds. \quad (9)$$

This equation looks similar to the stochastic case (4) although there are two main differences. The first is that the system is deterministic. The second is that it is defined for action functions and not for policies. We also define the optimal value of the deterministic limit $v_*(m_0)$:

$$v_*(m_0) = \sup_{\alpha} v_\alpha(m_0), \quad (10)$$

where the supremum is taken over all possible action functions from $[0; T] \rightarrow \mathcal{A}$.

2.5 Table of Notations

We recall here a list of the main notations used throughout the paper.

$M_\pi^N(k)$ Empirical measure of the system with N objects, under π , at time k , (Section 2.2)
$F^N(m, a)$ Drift of the system with N objects when the state is m and the action is a , Eq.(6)
$f(m, a)$ Drift of the limiting system (limit of rescaled $F^N(m, a)$ as $N \rightarrow \infty$), Eq.(11)
$\Phi_t(m_0, \alpha)$ State of the limiting system: $\Phi_t(m_0, \alpha) = m_0 + \int_0^t f(\Phi_s(m_0, \alpha), \alpha(s)) ds$, Eq.(8)
π^N Policy for the system with N objects: associates an action $a \in \mathcal{A}$ to each k , $M^N(k)$
α Action function for the limiting system: associates an action to each t : $\alpha : [0; T] \rightarrow \mathcal{A}$
π_*^N Optimal policy for the system with N objects
α_* Optimal action function for the limiting system (if it exists)
$V_\pi^N(m)$.. Expected reward for the system with N objects starting from m under policy π , Eq.(4)
$V_*^N(m)$ Optimal expected value for the system N : $V_*^N(m) = \sup_{\pi} V_\pi^N(m) = V_{\pi_*^N}^N(m)$, Eq.(5)
$V_\alpha^N(m)$ Expected value for the system N when applying the action function α , Eq.(12)
$v_\alpha(m)$ Value of the limiting system starting from m under action function α , Eq.(9)
$v_*(m)$ Optimal value of the limiting system: $v_*(m) = \sup_{\alpha} v_\alpha(m) = v_{\alpha_*}(m)$, Eq.(10)

2.6 Summary of Assumptions

In Section 3 we establish theorems for the convergence of the discrete stochastic optimization problem to a continuous deterministic one. These theorems are based on several technical assumptions, which are given next. Since SS is finite, the set $\mathcal{P}(SS)$ is the simplex in \mathbb{R}^{SS} and for $m, m' \in \mathcal{P}(SS)$ we define $\|m\|$ as the ℓ^2 -norm of m and $\langle m, m' \rangle = \sum_{i=1}^S m_i m'_i$ as the usual inner product.

(A1) (Transition probabilities) Objects can be observed only through their state, i.e., the transition probability matrix (or transition kernel) Γ^N , defined by Eq.(2), is invariant under permutations of $1 \dots N$.

There exist some non-random functions $I_1(N)$ and $I_2(N)$ such that $\lim_{N \rightarrow \infty} I_1(N) = \lim_{N \rightarrow \infty} I_2(N) = 0$ and such that for all m and any policy π , the number of objects that perform a transition between time slot k and $k + 1$ per time slot $\Delta_\pi^N(k)$ satisfies

$$\begin{aligned} \mathbb{E}(\Delta_\pi^N(k) | M_\pi^N(k) = m) &\leq N I_1(N) \\ \mathbb{E}(\Delta_\pi^N(k)^2 | M_\pi^N(k) = m) &\leq N^2 I(N) I_2(N) \end{aligned}$$

where $I(N)$ is the intensity function of the model, defined in the following assumption A2.

(A2) (Convergence of the Drift) There exist some non-random functions $I(N)$ and $I_0(N)$ and a function $f(m, a)$ such that $\lim_{N \rightarrow \infty} I(N) = \lim_{N \rightarrow \infty} I_0(N) = 0$ and

$$\left\| \frac{1}{I(N)} F^N(m, a) - f(m, a) \right\| \leq I_0(N) \quad (11)$$

f is defined on $\mathcal{P}(SS) \times \mathcal{A}$ and there exists L_2 such that $|f(m, a)| \leq L_2$.

(A3) (Lipschitz Continuity) There exist constants L_1 , K and K_r such that for all $m, m' \in \mathcal{P}(SS)$, $a, a' \in \mathcal{A}$:

$$\begin{aligned} \|F^N(m, a) - F^N(m', a)\| &\leq L_1 \|m - m'\| I(N) \\ \|f(m, a) - f(m', a')\| &\leq K(\|m - m'\| + d(a, a')) \\ |r(m, a) - r(m', a)| &\leq K_r \|m - m'\| \end{aligned}$$

We also assume that the reward is bounded: $\sup_{m, a \in \mathcal{A}} |r(m, a)| \stackrel{\text{def}}{=} \|r\|_\infty < \infty$.

To make things more concrete, here is a simple but useful case where all assumptions are true.

- There are constants c_1 and c_2 such that the expectation of the number of objects that perform a transition in one time slot is $\leq c_1$ and its standard deviation is $\leq c_2$,
- and $F^N(m, a)$ can be written under the form $\frac{1}{N} \varphi(m, a, 1/N)$ where φ is a continuous function on $\Delta_S \times \mathcal{A} \times [0, \epsilon)$ for some neighborhood Δ_S of $\mathcal{P}(SS)$ and some $\epsilon > 0$, continuously differentiable with respect to m .

In this case we can choose $I(N) = 1/N$, $I_0(N) = c_0/N$ (where c_0 is an upper bound to the norm of the differential $\frac{\partial \varphi}{\partial m}$), $I_1(N) = c_1/N$ and $I_2(N) = (c_1^2 + c_2^2)/N$.

3 Mean Field Convergence

In Section 3.1 we establish the main results, then, in Section 3.2, we provide the details of the method used to derive them.

3.1 Main Results

The first result establishes convergence of the optimization problem for the system with N objects to the optimization problem of the mean field limit:

Theorem 2 (Optimal System Convergence). *Assume (A0) to (A3). If $\lim_{N \rightarrow \infty} M^N(0) = m_0$ almost surely [resp. in probability] then:*

$$\lim_{N \rightarrow \infty} V_*^N(M^N(0)) = v_*(m_0)$$

almost surely [resp. in probability], where V_^N and v_* are the optimal values for the system with N objects and the mean field limit, defined in Section 2.*

The proof is given in Section 5.6.

The second result states that an optimal action function for the mean field limit provides an asymptotically optimal strategy for the system with N objects. We need, at this point, to introduce a first auxiliary system, which is a system with N objects controlled by an action function borrowed from the mean field limit. More precisely, let α be an action function that specifies the action to be taken at time t . Although α has been defined for the limiting system, it can also be used in the system with N objects. In this case, the action function α can be seen as a policy that does not depend on the state of the system. At step k , the controller applies action $\alpha(kI(N))$. By abuse of notation, we denote by M_α^N , the state of the system when applying the action function α (it will

be clear from the notation whether the subscript is an action function or a policy). The value for this system is defined by

$$V_\alpha^N(m_0) \stackrel{\text{def}}{=} \mathbb{E} \left(\sum_{k=0}^{H^N} r(M_\alpha^N(k), \alpha(kI(N))) \middle| M_\alpha^N(0) = m_0 \right) \quad (12)$$

Our next result is the convergence of convergence of M_α^N and of the value:

Theorem 3. *Assume (A0) to (A3); α is a piecewise Lipschitz continuous action function on $[0; T]$, of constant K_α , and with at most p discontinuity points. Let $\hat{M}_\alpha^N(t)$ be the linear interpolation of the discrete time process M_α^N . Then for all $\epsilon > 0$:*

$$\mathbb{P} \left\{ \sup_{0 \leq t \leq T} \left\| \hat{M}_\alpha^N(t) - \phi_t(m_0, \alpha) \right\| > [\|M^N(0) - m_0\| + I'_0(N, \alpha)T + \epsilon] e^{L_1 T} \right\} \leq \frac{J(N, T)}{\epsilon^2} \quad (13)$$

and

$$|V_\alpha^N(M^N(0)) - v_\alpha(m_0)| \leq B'(N, \|M^N(0) - m_0\|) \quad (14)$$

where J, I'_0 and B' are defined in Section 5.1 and satisfy $\lim_{N \rightarrow \infty} I'_0(N, \alpha) = \lim_{N \rightarrow \infty} J(N, T) = 0$ and $\lim_{N \rightarrow \infty, \delta \rightarrow 0} B'(N, \delta) = 0$.

In particular, if $\lim_{N \rightarrow \infty} M_\pi^N(0) = m_0$ almost surely [resp. in probability] then $\lim_{N \rightarrow \infty} V_\alpha^N(M^N(0)) = v_\alpha(m_0)$ almost surely [resp. in probability].

The proof is given in Section 5.5.

As the reward function $r(m, a)$ is bounded and the time-horizon $[0; T]$ is finite, the set of values when starting from the initial condition m , $\{v_\alpha(m) : \alpha \text{ action function}\}$, is bounded. This set is not necessarily compact because the set of action functions may not be closed (a limit of Lipschitz continuous functions is not necessarily Lipschitz continuous). However, as it is bounded, for all $\epsilon > 0$, there exists an action function α^ϵ such that $v_*(m) = \sup_\alpha v_\alpha(m) \leq v_{\alpha^\epsilon} + \epsilon$. Theorem 2 shows that α^ϵ is optimal up to 2ϵ for N large enough. This shows the following corollary:

Corollary 4 (Asymptotically Optimal Policy). *Let α^* be an optimal action function for the limiting system. Then*

$$\lim_{N \rightarrow \infty} |V_{\alpha^*}^N - V_*^N| = 0$$

In other words, an optimal action function for the limiting system is asymptotically optimal for the system with N objects.

In particular, this shows that as N grows, policies that do not take into account the state of the system (*i.e.*, action functions) are asymptotically as good as adaptive policies. In practice however, adaptive policies might perform better, especially for very small values of N . However, it is in general impossible to prove convergence for adaptive policies.

3.2 Derivation of Main Results

3.2.1 Second Auxiliary System

The method of proof uses a second auxiliary system, the process $\phi_t(m_0, A_\pi^N)$ defined below. It is a limiting system controlled by an action function derived from the policy of the original system with N objects.

Consider the system with N objects under policy π . The process M_π^N is defined on some probability space Ω . To each $\omega \in \Omega$ corresponds a trajectory $M_\pi^N(\omega)$, and for each $\omega \in \Omega$, we define an action function $A_\pi^N(\omega)$. This random function is piecewise constant on each interval $[kI(N), (k+1)I(N))$ ($k \in \mathbb{N}$) and is such that $A_\pi^N(\omega)(kI(N)) \stackrel{\text{def}}{=} \pi_k(M^N(k))$ is the action taken by the controller of the system with N objects at time slot k , under policy π .

Recall that for any $m_0 \in \mathcal{P}(SS)$ and any action function α , $\phi_t(m_0, \alpha)$ is the solution of the ODE (7). For every ω , $\phi_t(m_0, A_\pi^N(\omega))$ is the solution of the limiting system with action function $A_\pi^N(\omega)$, i.e.

$$\phi_t(m_0, A_\pi^N(\omega)) - m_0 = \int_0^t f(\phi_s(m_0, A_\pi^N(\omega)), A_\pi^N(\omega)(s)) ds.$$

When ω is fixed, $\phi_t(m_0, A_\pi^N(\omega))$ is a continuous time deterministic process corresponding to one trajectory $M_\pi^N(\omega)$. When considering all possible realizations of M_π^N , $\phi_t(m_0, A_\pi^N)$ is a random, continuous time function “coupled” to M_π^N . Its randomness comes only from the action term A_π^N , in the ODE. In the following, we omit to write the dependence in ω . A_π^N and M_π^N will always designate the processes corresponding to the same ω .

3.2.2 Convergence of Controlled System

The following result is the main technical result; it shows the convergence of the controlled system in probability, with explicit bounds. Notice that it does not require any regularity assumption on the policy π .

Theorem 5. *Under Assumptions (A0) to (A3), for any $\epsilon > 0$, $N \geq 1$ and any policy π :*

$$\mathbb{P} \left\{ \sup_{0 \leq t \leq T} \left\| \hat{M}_\pi^N(t) - \phi_t(m_0, A_\pi^N) \right\| > [\|M^N(0) - m_0\| + I_0(N)T + \epsilon] e^{L_1 T} \right\} \leq \frac{J(N, T)}{\epsilon^2} \quad (15)$$

where \hat{M}_π^N is the linear interpolation of the discrete time system with N objects) and J is defined in Section 5.1.

Recall that $I_0(N)$ and $J(N, T)$ for a fixed T go to 0 as $N \rightarrow \infty$. The proof is given in Section 5.3.

3.2.3 Convergence of Value

Let π be a policy and A_π^N the sequence of actions corresponding to a trajectory M_π^N as we just defined. Eq.(9) defines the value for the deterministic limit when applying a sequence of actions. This defines a random variable $v_{A_\pi^N}(m_0)$ that corresponds to the value over the limit system when using A_π^N as action function. The random part comes from A_π^N . $\mathbb{E}[v_{A_\pi^N}(m_0)]$ designates the expectation of this value over all possible A_π^N . A first consequence of Theorem 5 is the convergence of $V_\pi^N(M^N(0))$ to $\mathbb{E}[v_{A_\pi^N}(m_0)]$ with an error that can be uniformly bounded.

Theorem 6 (Uniform convergence of the value). *Let A_π^N be the random action function associated with M_π^N , as defined earlier. Under Assumptions (A0) to (A3),*

$$|V_\pi^N(M^N(0)) - \mathbb{E}[v_{A_\pi^N}(m_0)]| \leq B(N, \|M^N(0) - m_0\|)$$

where B is defined in Section 5.1.

Note that $\lim_{N \rightarrow \infty, \delta \rightarrow 0} B(N, \delta) = 0$; in particular, if $\lim_{N \rightarrow \infty} M_\pi^N(0) = m_0$ almost surely [resp. in probability] then $|V_\pi^N(M^N(0)) - \mathbb{E}[v_{A_\pi^N}(m_0)]| \rightarrow 0$ almost surely [resp. in probability].

The proof is given in Section 5.4.

3.2.4 Putting Things Together

The proof of the main result uses the two auxiliary systems. The first auxiliary system provides a strategy for the system with N objects derived from an action function of the mean field limit; it cannot do better than the optimal value for the system with N objects, and is close to the optimal value of the mean field limit. Therefore, the optimal value for the system with N objects is lower bounded by the optimal value for the mean field limit. The second auxiliary system is used in the opposite direction, which shows that, roughly speaking, for large N the two optimal values are the same. We give the details of the derivation in Section 5.6.

4 Applications

4.1 Hamilton-Jacobi-Bellman Equation and Dynamic Programming

Let us now consider the finite time optimization problem for the stochastic system and its limit from a constructive point of view. As the state space is finite, we can compute the optimal value by using a dynamic programming algorithm. If $U^N(m, t)$ denotes the optimal value for the stochastic system starting from m at time $t/I(N)$, then $U^N(m, t) = \sup_{\pi} \mathbb{E} \left[\sum_{k=t/I(N)}^{T/I(N)} r^N(M_{\pi}^N(k)) : M^N(t) = m \right]$. The optimal value can be computed by a discrete dynamic programming algorithm [21] by setting $U^N(m, T) = r^N(m)$ and

$$U^N(m, t) = \sup_{a \in \mathcal{A}} \mathbb{E} \left(r^N(m, a) + U^N(M^N(t + I(N)), t + I(N)) \mid \bar{M}^N(t) = m, A^N(t) = a \right). \quad (16)$$

Then, the optimal cost over horizon $[0; T/I(N)]$ is $V_*^N(m) = U(m, 0)$.

Similarly, if we denote by $u(m, t)$ the optimal cost over horizon $[t; T]$ for the limiting system, $u(m, t)$ satisfies the classical Hamilton-Jacobi-Bellman equation:

$$\dot{u}(m, t) + \max_a \{ \nabla u(m, t) \cdot f(m, a) + r(m, a) \} = 0. \quad (17)$$

This provides a way to compute the optimal value, as well as the optimal policy, by solving the partial differential equation above.

4.2 Algorithms

Theorem 2 above can be used to design an effective construction of an asymptotically optimal policy for the system with N objects over the horizon $[0, H]$ by using the procedure described in Algorithm 1.

Algorithm 1: Static algorithm constructing a policy for the system with N objects, over the finite horizon.

begin

From the original system with N objects, construct the occupancy measure M^N and its kernel Γ^N and let $M^N(0)$ be the initial occupancy measure;

Compute the limit of the drift of Γ^N , namely the function f ;

Solve the HJB equation (17) on the interval $[0, HI(N)]$. This provides an optimal control function $\alpha(M_0^N, t)$;

Construct a discrete control $\pi(M^N(k), k)$ for the discrete system, that gives the action to be taken under state $M^N(k)$ at step k :

$$\pi(M^N(k), k) \stackrel{\text{def}}{=} \alpha(\phi_{kI(N)}(M^N(0), \alpha)).$$

return π ;

Theorem 2 says that under policy π , the total value V_{π}^N is asymptotically optimal:

$$\lim_{N \rightarrow \infty} V_{\pi}^N(M^N(0)) = \liminf_{N \rightarrow \infty} V_*^N(M^N(0)).$$

The policy π constructed by Algorithm 1 is static in the sense that it does not depend on the state $M^N(k)$ but only on the initial state $M^N(0)$, and the deterministic estimation of $M^N(k)$ provided by the differential equation. One can construct a more adaptive policy by updating the starting point of the differential equation at each step. This new procedure, constructing an adaptive policy π' from 0 to the final horizon H is given in Algorithm 2.

In practice, the total value of the adaptive policy π' is larger than the value of the static policy π because it uses on-line corrections at each step, before taking a new action. However Theorem 2 does not provide a proof of its asymptotic optimality.

Algorithm 2: Adaptive algorithm constructing a policy π' for the system with N objects, over the finite horizon H .

```

begin
   $M := M^N(0); k := 0;$ 
  repeat
     $\alpha_k(M, \cdot) :=$  solution of (17) over  $[kI(N), HI(N)]$  starting in  $M$ ;
     $\pi'(M, k) := \alpha_k(\phi_{kI(N)}(M, \alpha_k));$ 
     $M$  is changed by applying kernel  $\Gamma_{\pi'}^N$ ;
     $k := k+1;$ 
  until  $k=H;$ 
  return  $\pi';$ 

```

4.3 Examples

In this section, we develop three examples. The first one can be seen as a simple illustration of optimal mean field. The limiting ODE is quite simple and can be optimized in closed analytical form.

The second example considers a classic virus problem. Although virus propagations concern discrete objects (individuals or devices), most work in the literature study a continuous approximation of the problem under the form of an ODE. The justification of passing from a discrete to a continuous model is barely mentioned in most papers (they mainly focus on the study of the ODE). Here we present a discrete dynamical system based on a simple stochastic mobility model for the individuals whose behavior converges to a classic continuous model. We show on a numerical example that the limiting problem provides a policy that is close to optimal, even for a system with a relatively small numbers of nodes.

Finally, the last example comes from routing optimization in a queueing network model of volunteer computing platforms. The purpose of this last example is to show that a discrete optimal control problem suffering from the curse of dimensionality can be replaced by a continuous optimization problem where an HJB equation must be solved over a much smaller state space.

4.3.1 Utility Provider Pricing

This is a simplified discrete Merton's problem. This example shows a case where the optimization problem in the infinite system can be solved in closed form. This can be seen as an ideal case for the mean field approach: although the original system is difficult to solve even numerically when N is large, taking the limit when N goes to infinity makes it simple to solve, in an analytical form.

We consider a system made of a utility and N users; users can be either in state S (subscribed) or U (unsubscribed). The utility fixes their price $\alpha \in [0, 1]$. At every time step, one randomly chosen customer revises her status: if she is in state U [resp. S], with probability $s(\alpha)$ [resp. $a(\alpha)$] she moves to the other state; $s(\alpha)$ is the probability of a new subscription, and $a(\alpha)$ is the probability of attrition. We assume $s(\cdot)$ decreases with α and $a(\cdot)$ increases. If the price is large, the instant gain is large, but the utility loses customers, which eventually reduces the gain.

Within our framework, this problem can be seen as a Markovian system made of N objects (users) and one controller (the provider). The intensity of the model is $I(N) = 1/N$. Moreover, if the immediate profit is divided by N (this does not alter the optimal pricing policy) and if $x(t)$ is the fraction of objects in state S at time t and $\alpha(t) \in [0, 1]$ is the action taken by the provider at time t , the mean field limit of the system is:

$$\frac{\partial x}{\partial t} = -x(t)a(\alpha(t)) + (1 - x(t))s(\alpha(t)) = s(\alpha(t)) - x(s(\alpha(t)) + a(\alpha(t))) \quad (18)$$

and the rescaled profit over a time horizon T is $\int_0^T x(t)\alpha(t)dt$. Call $u_*(t, x)$ the optimal benefit over the interval $[t, T]$ if there is a proportion x of subscribers at time t . The Hamilton-Jacobi-Bellman

equation is

$$\begin{aligned} \frac{\partial}{\partial t} u_*(t, x) + H\left(x, \frac{\partial}{\partial x} u_*(t, x)\right) &= 0 \\ \text{with } H(x, p) &= \max_{\alpha \in [0,1]} [p(s(\alpha) - x(s(\alpha) + a(\alpha)) + \alpha x] \end{aligned}$$

H can be computed under reasonable assumptions on the rates of subscription and attrition $s()$ and $a()$, which can then be used to show that there exists an optimal policy that is threshold based. To continue the rest of this illustration, we consider the radically simplified case where α can take only the values 0 and 1 and under the conditions $s(0) = a(1) = 1$ and $s(1) = a(0) = 0$, in which case the ODE becomes

$$\frac{\partial x}{\partial t} = -x(t)\alpha(t) + (1 - x(t))(1 - \alpha(t)) = 1 - x(t) - \alpha(t), \quad (19)$$

and $H(x, p) = \max(x(1 - p), (1 - x)p)$. The solution of the HJB equation can be given in closed form. The optimal policy is to choose action $\alpha = 1$ if $x > 1/2$ or $x > 1 - \exp(-(T - t))$, and 0 otherwise. Figure 1 shows the evolution of the proportion of subscribers $x(t)$ when the optimal policy is used. The coloured area corresponds to all the points (t, x) where the optimal policy is $\alpha = 1$ (fix a high price) and the white area is where the optimal policy is to choose $\alpha = 0$ (low price).

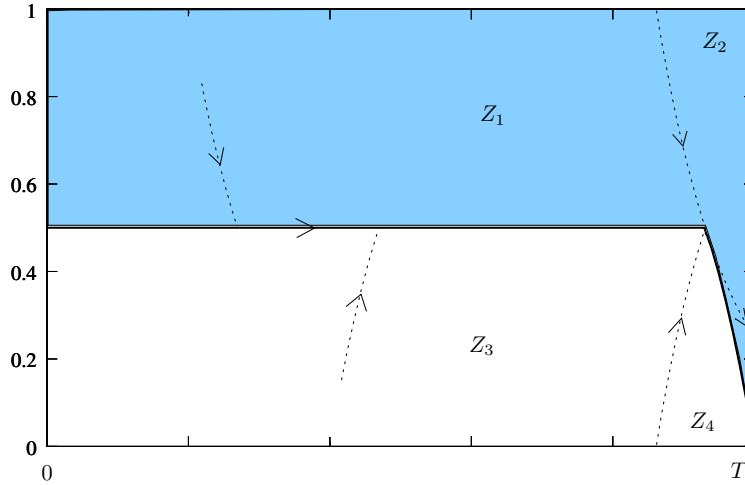


Figure 1: Evolution of the proportion of subscribers (y -axis) under the optimal pricing policy.

To show that this policy is indeed optimal, one has to compute the corresponding value of the benefit $u(t, x)$ and show that it satisfies the HJB equation. This can be done using a case analysis, by computing explicitly the value of $u(t, x)$ in the zones Z_1, Z_2, Z_3 and Z_4 displayed in Figure 1, and check that $u(t, x)$ satisfies Eq.(19) in each case.

4.3.2 Infection Strategy of a Viral Worm

This second example has two purposes. The first one is to provide a rigorous justification of the use of a continuous optimization approach for this classic problem in population dynamics and to show that the continuous limit provides insights on the structure of the optimal behavior for the discrete system. Here, the optimal action function can be shown to be of the bang-bang type for the limit problem, by using tools from continuous optimization such as the Pontryagin maximum principle. Theorem 2 shows that a bang-bang policy should also be asymptotically optimal in the discrete case.

The second purpose is to compare numerically the performance of the optimal policy of the deterministic limit α_* and the performance of other policies for the stochastic system for small values of N . We show that α_* is close to optimal even for $N = 10$ and that it outperforms another classic heuristic.

This example is taken from [15] and considers the propagation of infection by a viral worm. Actually, similar epidemic models have been validated through experiments, as well as simulations as a realistic representation of the spread of a virus in mobile wireless networks (see [7, 22]). A *susceptible* node is a mobile wireless device, not contaminated by the worm but prone to infection. A node is *infective* if it is contaminated by the worm. An infective node spreads the worm to a susceptible node whenever they meet, with probability β . The worm can also choose to kill an infective node, i.e., render it completely dysfunctional - such nodes are denoted *dead*. A functional node that is immune to the worm is referred to as *recovered*. Although the network operator uses security patches to immunize susceptibles (they become recovered) and heals infectives to the recovered state, the goal of the worm is to maximize the damages done to the network. Let the total number of nodes in the network be N . Let the proportion of susceptible, infective, recovered and dead nodes at time t be denoted by $S(t)$, $I(t)$, $R(t)$ and $D(t)$, respectively. Under a uniform mobility model, the probability that a susceptible node becomes infected is $\beta I/N$. The immunization of susceptibles (resp. infectives) happens at a fixed rate q (resp. b). This means that a susceptible (resp. infective) node is immunized with probability q/N (resp. b/N) at every time step.

At this point, authors of [15] invoke the classic results of Kurtz [17] to show that the dynamics of this population process converges to the solution of the following differential equations.

$$\begin{aligned} \frac{\partial S}{\partial t} &= -\beta IS - qS \\ \frac{\partial I}{\partial t} &= \beta IS - bI - v(t)I \\ \frac{\partial D}{\partial t} &= v(t)I \\ \frac{\partial R}{\partial t} &= bI + qS. \end{aligned} \tag{20}$$

This system actually satisfies assumptions (A_1, A_2, A_3) , which allows us not only to obtain the mean field limit, but also to say more about the optimization problem. The objective of the worm is to find $v(\cdot)$ such that the damage function $D(T) + \int_0^T f(I(t))dt$ is maximized under the constraint $0 \leq v \leq v_{\max}$ (where f is convex). In [15], this problem is shown to have a solution and the Pontryagin maximum principle is used to show that the optimal solution $v_*(\cdot)$ is of bang-bang type:

$$\exists t_1 \in [0 \dots T] \text{ s. t. } v_*(t) = 0 \text{ for } 0 < t < t_1 \text{ and } v_*(t) = v_{\max} \text{ for } t_1 < t < T. \tag{21}$$

Theorem 2 makes the formal link between the optimization of the model on an individual level and the previous resolution of the optimization problem on the differential equations, done in [15]. It allows us to formally claim that the policy α_* of the worm is indeed asymptotically optimal when the number of objects goes to infinity.

We investigated numerically the performance of α_* against various infection policies for small values of the number of nodes in the system N . These results are reported on Figure 2, where we compare four values:

- v_* – the optimal value of the limiting system.
- V_*^N – the optimal expected damage for the system with N objects (MDP problem);
- $V_{\alpha_*}^N$ – the expected value of the system with N objects when applying the action function α_* that is optimal for the limiting system; Performance of algorithm 1
- the performance of a heuristic where, instead of choosing a threshold as suggested by the limiting system (21), the killing probability ν is fixed for the whole time. The curve on the figure is drawn for the optimal ν (recomputed for each parameter N).

We implemented a simulator that follows strictly the model of infection described earlier in this part. We chose parameters similar to those used in [15]: the parameter for the evolution of the system are

$\beta = .6$, $q = .1$, $b = .1$, $v_{\max} = 1$ and the damage function to be optimized is $D(T) + \frac{1}{T} \int_0^T I^2(t)dt$ with $T = 10$. However, it should be noted that the choice of these parameters does not influence qualitatively the results. Thanks to the relatively small size of the system, these four quantities can be computed numerically using a backward induction. The optimal policies for the deterministic limit consists in not killing machines until $t_1 = 4.9$ and in killing machines at a maximum rate after that time: $\alpha_*(t) = \mathbf{1}_{\{t > 4.9\}}$.

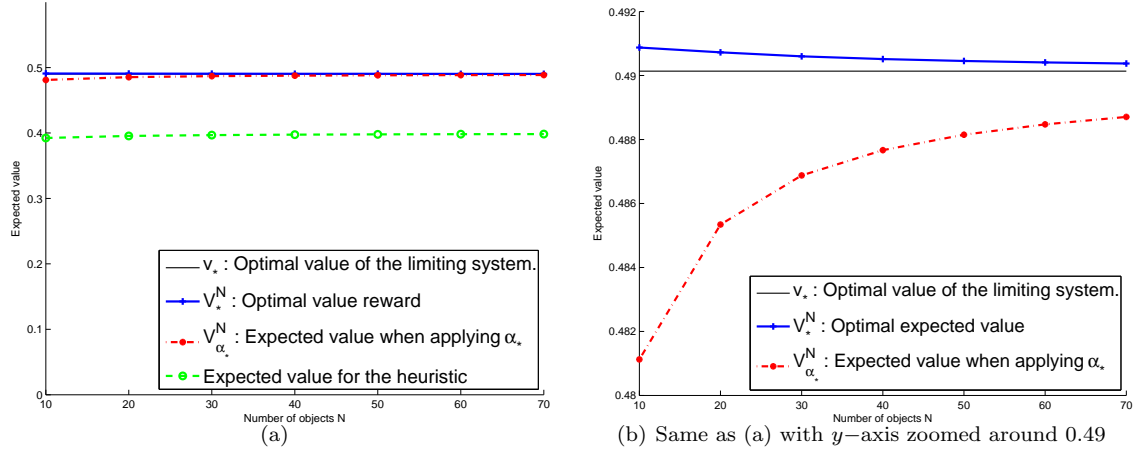


Figure 2: Damage caused by the worm for various infection policies as a function of the size of the system N . The goal of the worm is to *maximize* the damage (higher means better). Panel (a) shows the optimal value v_* for the limiting system (mean field limit), the optimal value V_*^N for the system with N objects, the value $v_{\alpha_*}^N$ of the asymptotically optimal policy given in Corollary 4 and the value of a classic heuristic. Panel (b) zooms the y -axis around the values of the optimal policies.

Theorem 2 shows that α_* is asymptotically optimal ($\lim_{N \rightarrow \infty} V_*^N = \lim_{N \rightarrow \infty} V_{\alpha_*}^N = v_*$), but Figure 2(a) shows that, already for low values of N , these three quantities are very close. A classic heuristic for this maximal infection problem is to kill a node with a constant probability ν , regardless of the time horizon. Our numerical study shows that α_* outperforms this heuristic by more than 20%. The performance of this heuristic does not increase with the size of the system N .

In order to illustrate the convergence of the values V_*^N and $V_{\alpha_*}^N$ to v_* , Figure 2(b) is a detailed view of Figure 2(a) where we show the two quantities V_*^N , $V_{\alpha_*}^N$ and their common limit v_* . This figure shows that the convergence is indeed very fast. Other numerical experiments indicate that this is true for a large panel of parameters. Although this figure seems to indicate that $V_{\alpha_*}^N \leq v_* \leq V_*^N$, this is not true in general, for example adding $5D(t)$ to the damage function leads to $V_{\alpha_*}^N \leq V_*^N \leq v_*$ ($V_{\alpha_*}^N$ is always less than V_*^N by definition of V_*^N).

4.3.3 Brokering Problem

Finally, let us consider a model of a volunteer computing system such as BOINC <http://boinc.berkeley.edu/>. Volunteer computing means that people make their personal computer available for a computing system. When they do not use their computer, it is available for the computing system. However, as soon as they start using their computer, it becomes unavailable for the computing system. These systems are becoming more and more popular and provide large computing power at a very low cost [16].

The Markovian model with N objects is defined as follows. The N objects represent the users that can submit jobs to the system and the resources that can run the jobs. The resources are grouped into a small number of clusters and all resources in the same cluster share the same characteristics in terms of speed and availability. Users send jobs to a central broker whose role is to balance the load among the clusters.

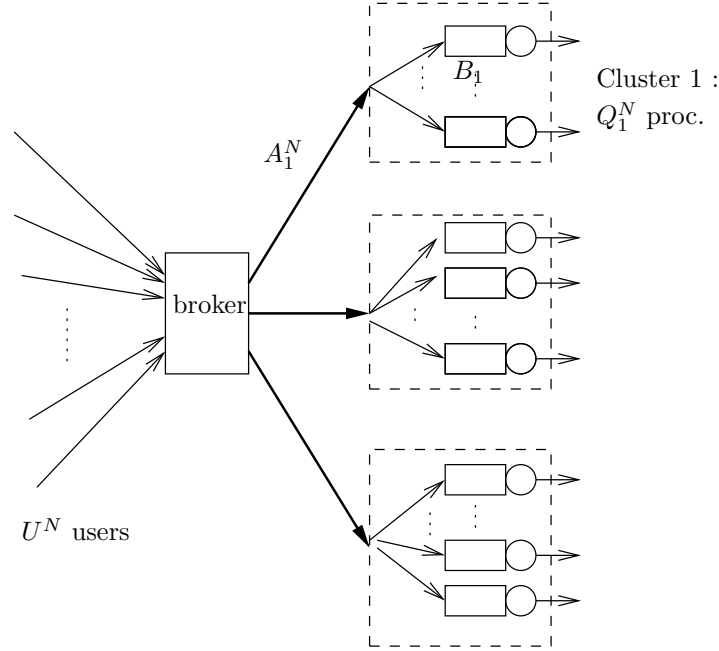


Figure 3: The brokering problem in a desktop grid system, such as Boinc

The model is a discrete time model of a queuing system. Actually, a more natural continuous-time Markov model could also be handled similarly, by using uniformization.

There are U^N users. Each user has a state $x \in \{\text{on}, \text{off}\}$. At each time step, an active user sends one job with probability p_s^N and becomes inactive with probability p_i/N . An inactive user sends no jobs to the system and becomes **on** with probability p_o/N .

There are C clusters in the system. Each cluster c contains Q_c^N computing resources. Each resource has a buffer of bounded size J_c . A resource can either be valid or broken. If it is valid and if it has one or more job in its queue, it completes one job with probability μ_c/N at this time slot. A resource gets **broken** with probability p_b/N . In that case, it discards all the packets of its buffer. A broken resource becomes **valid** with probability p_v/N .

At each time step, the broker takes an action $a \in \mathcal{P}(\{1 \dots C\})$ and sends the packets it received to the clusters according to the distribution a . A packet sent to cluster c joins the queue of one of the resources, say k ; according to a local rule (for example chosen uniformly among the Q_c^N resources composing the cluster). If the queue of resource k is full, the packet is lost. The goal of the broker is to minimize the number of losses plus the total size of the queues over a finite horizon (and hence the response time of accepted packets).

This model is represented in Figure 3.

The system has an intensity $I(N) \stackrel{\text{def}}{=} 1/N$. The number C of clusters is fixed and does not depend on N , as well as the sizes J_c of the buffers. However, both the number of users U^N , and the number of resources in the clusters Q_c^N , are linear in N . Finally, by construction, all the state changes occur with probabilities that scale with $1/N$.

The limiting system is described by the variable $m_o(t)$, that represents the fraction of users who are on, and the variables $q_{c,j}(t)$ and $b_c(t)$ that, respectively, represent the fraction of resources in cluster c having j jobs in their buffer and the fraction of resources in cluster c that are broken. For an action function $\alpha(\cdot)$, we denote by $\alpha_c(\cdot)$ the fraction of packets sent to cluster c . Finally, let us denote by m the fraction of users (both active or inactive) and q_c the fraction of processors in cluster c . These fractions are constant (independent of time) and satisfy $m + q_1 + \dots + q_C = 1$.

We get the following equations:

$$\frac{\partial m_o(t)}{\partial t} = -p_i m_o(t) + p_o(m - m_o(t)) \quad (22)$$

$$\frac{\partial q_{c,0}(t)}{\partial t} = p_a b_c(t) - \frac{\alpha_c(t) p_s m_o(t)}{q_c} q_{c,0}(t) + \mu_c q_{c,1} - p_b q_{c,0}(t) \quad (23)$$

$$\frac{\partial q_{c,j}(t)}{\partial t} = \frac{\alpha_c(t) p_s m_o(t)}{q_c} (q_{c,j-1}(t) - q_{c,i}(t)) + \mu_c (q_{c,j+1} - q_{c,j}) - p_b q_{c,j}(t) \quad (24)$$

$$\frac{\partial q_{c,J_c}(t)}{\partial t} = \frac{\alpha_c(t) p_s m_a(t)}{q_c} q_{c,J_c-1}(t) - \mu_c q_{c,J_c} - p_b q_{c,J_c}(t) \quad (25)$$

$$\frac{\partial b_c(t)}{\partial t} = -p_v b_c(t) + p_b \sum_{j=0}^{J_c} q_{c,j}(t). \quad (26)$$

where (23) and (25) hold for each cluster c and (24) holds for each cluster c and for all $j \leq J_c$. The cost associated to the action function α is:

$$\int_0^T \sum_{c=1}^C \sum_{j=1}^{J_c} j q_{c,j}(t) + \gamma \left(\sum_{c=1}^C \frac{\alpha_c(t) p_s m_o(t)}{q_c} (q_{c,J_c}(t) + b_c(t)) + \sum_{c=1}^C p_b \sum_{j=1}^{J_c} j q_{c,j}(t) \right) dt \quad (27)$$

The first part of (27) represents the cost induced by the number of jobs in the system. The second part of (27) represents the cost induced by the losses. The parameter γ gives weight on the cost induced by the losses.

The HJB problem becomes minimizing (27) subject to the variables $u_a, q_{k,i}, b_k$ satisfying Equations (22) to (26). This system is made of $(J+2)C$ ODEs. Solving the HJB equation numerically in this case can be challenging but remains more tractable than solving the original Bellman equation over J^N states. The curse of dimensionality is so acute for the discrete system that it cannot be solved numerically with more than 10 processors [5].

5 Proofs

5.1 Details of Scaling Constants

$$\begin{aligned} I'_0(N, \alpha) &\stackrel{\text{def}}{=} I_0(N) + I(N) K e^{(K-L_1)T} \left(\frac{K\alpha}{2} + 2(1 + \min(1/I(N), p)) \|\alpha\|_\infty \right) \\ J(N, T) &\stackrel{\text{def}}{=} 8T \left\{ L_1^2 [I_2(N) I(N)^2 + I_1(N)^2 (T + I(N))] + S^2 [2I_2(N) + I(N) (I_0(N) + L_2)^2] \right\} \\ B(N, \delta) &\stackrel{\text{def}}{=} I(N) \|r\|_\infty + K_r (\delta + I_0(N)T) \frac{e^{L_1 T} - 1}{L_1} \\ &\quad + \frac{3}{2^{\frac{1}{3}}} \left[\frac{K_r}{L_1} \left(e^{L_1 T} - 1 + \frac{I(N)}{2} \right) \right]^{\frac{2}{3}} \|r\|_\infty^{\frac{1}{3}} J(N, T)^{\frac{1}{3}} \\ B'(N, \delta) &\stackrel{\text{def}}{=} I(N) \|r\|_\infty + K_r [\delta + I'_0(N, \alpha)T] \frac{e^{L_1 T} - 1}{L_1} \\ &\quad + \frac{3}{2^{\frac{1}{3}}} \left[\frac{K_r}{L_1} \left(e^{L_1 T} - 1 + \frac{I(N)}{2} \right) \right]^{\frac{2}{3}} \|r\|_\infty^{\frac{1}{3}} J(N, T)^{\frac{1}{3}} \end{aligned}$$

5.2 Proof of Theorem 1

We begin with a few general statements. Let \mathcal{P} be the set of probabilities on SS and $\mu^N : SS^N \rightarrow \mathcal{P}$ defined by $\mu^N(x)_i = \frac{1}{N} \sum_{n=1}^N \mathbf{1}_{x_n=i}$ for all $i \in SS$. Also let \mathcal{P}^N be the image set of μ^N , i.e the set of all occupancy measures that are possible when the number of objects is N . The following

establishes that if two global states have the same occupancy measure, then they differ by a permutation.

Lemma 1. *For all $x, x' \in SS^N$, if $\mu^N(x) = \mu^N(x')$ there exists some $\sigma \in \mathfrak{S}^N$ such that $x' = \sigma(x)$.*

Proof. By induction on N . Its is obvious for $N = 1$. Assume the lemma holds for $N - 1$ and let $x, x' \in SS^N$, with $\mu^N(x) = \mu^N(x')$. There is at least one coordinate, say i , such that $x'_i = x_1$, because there is the same number of occurrences of $s = x_1$ in both x and x' . Let $y = x_2 \dots x_N$ and $y' = x'_2 \dots x'_N$. Then $\mu^{N-1}(y) = \mu^{N-1}(y')$, therefore there exists some $\tau \in \mathfrak{S}^{N-1}$ such that $y' = \tau(y)$. Define σ by $\sigma(1) = i$, $\sigma(j) = \tau(j) + \mathbf{1}_{\tau(j) > i}$, for $j \geq 2$, so that $x' = \sigma(x)$. Clearly σ is a permutation of $\{1, \dots, N\}$. \square

Let $f : SS^N \rightarrow E$ where E is some arbitrary set. We say that f is invariant under \mathfrak{S}^N if $f \circ \sigma = f$ for all $\sigma \in \mathfrak{S}^N$. The following results states that if a function of the global state is invariant under permutations, it is a function of the occupancy measure.

Lemma 2. *If $f : SS^N \rightarrow E$ is invariant under \mathfrak{S} then there exists $\bar{f} : \mathcal{P}^N \rightarrow E$ such that $\bar{f} \circ \mu^N = f$.*

Proof. Define \bar{f} as follows. For every $m \in \mathcal{P}^N$ pick some arbitrary $x_0 \in (\mu^N)^{-1}(m)$ and let $\bar{f}(m) = f(x_0)$. Now let x , perhaps different from x_0 , such that $\mu^N(x) = m$. By Lemma 1, there exists some $\sigma \in \mathfrak{S}^N$ such that $x = \sigma(x_0)$ therefore $f(x) = f(x_0) = \bar{f}(\mu^N(x))$. This is true for every $m \in \mathcal{P}^N$ thus $f(x) = \bar{f}(\mu^N(x))$ for every $x \in SS^N$. \square

The sequence of actions a_k is given and N is fixed. We are thus given a time-inhomogeneous Markov chain X^N on SS^N , with transition kernel G_k , $k \in \mathbb{N}$, given by $G_k(x, y) = \Gamma^N(x, y, a_k)$, such that for any permutation $\sigma \in \mathfrak{S}^N$ and any states x, y we have

$$G_k(\sigma(x), \sigma(y)) = G_k(x, y) \quad (28)$$

Let $\mathcal{F}(k)$ be the σ -field generated by $X^N(s)$ for $s \leq k$ and $\mathcal{G}(k)$ be the σ -field generated by $M^N(s)$ for $s \leq k$. Note that because $M^N = \mu^N \circ X^N$, $\mathcal{G}(k) \subset \mathcal{F}(k)$.

Pick some arbitrary test function $\varphi : SS^N \rightarrow \mathbb{R}$ and fix some time $k \geq 1$; we will now compute $\mathbb{E}(\varphi(M^N(k)) | \mathcal{F}(k-1))$. Because M^N is a function of X^N and X^N is a Markov chain, $\mathbb{E}(\varphi(M^N(k)) | \mathcal{F}(k-1))$ is a function, say ψ , of $X^N(k-1)$. We have, for any fixed $\sigma \in \mathfrak{S}^N$:

$$\begin{aligned} \psi(x) &\stackrel{\text{def}}{=} \sum_{y \in SS^N} G_k(x, y) \varphi(\mu^N(y)) = \sum_{y \in SS^N} G_k(x, \sigma(y)) \varphi(\mu^N(\sigma(y))) \\ &= \sum_{y \in SS^N} G_k(x, \sigma(y)) \varphi(\mu^N(y)) \\ \psi(\sigma(x)) &= \sum_{y \in SS^N} G_k(\sigma(x), \sigma(y)) \varphi(\mu^N(y)) = \sum_{y \in SS^N} G_k(x, y) \varphi(\mu^N(y)) \end{aligned}$$

where the last equality is by Eq.(28). Thus $\psi(\sigma(x)) = \psi(x)$ and by Lemma 2 there exists some function $\bar{\psi}$ such that $\psi(x) = \bar{\psi}(\mu^N(x))$, i.e.

$$\mathbb{E}(\varphi(M^N(k)) | \mathcal{F}(k-1)) = \bar{\psi}(M^N(k-1)) \quad (29)$$

In particular, $\mathbb{E}(\varphi(M^N(k)) | \mathcal{F}(k-1))$ is $\mathcal{G}(k-1)$ -measurable. Now

$$\begin{aligned} \mathbb{E}(\varphi(M^N(k)) | \mathcal{G}(k-1)) &= \mathbb{E}(\mathbb{E}(\varphi(M^N(k)) | \mathcal{F}(k-1)) | \mathcal{G}(k-1)) \\ &= \mathbb{E}(\bar{\psi}(M^N(k-1)) | \mathcal{G}(k-1)) = \bar{\psi}(M^N(k-1)) \end{aligned}$$

which expresses that M^N is a Markov chain.

5.3 Proof of Theorem 5

The proof is inspired by the method in [3]. The main idea of the proof is to write

$$\begin{aligned} \|M_\pi^N(k) - \phi_{kI(N)}(m_0, A_\pi^N)\| &\leq \left\| M_\pi^N(k) - M^N(0) - \sum_{j=0}^{k-1} f^N(j) \right\| \\ &\quad + \left\| M^N(0) + \sum_{j=0}^{k-1} f^N(j) - \phi_{kI(N)}(m_0, A_\pi^N) \right\| \end{aligned}$$

where $f^N(k) \stackrel{\text{def}}{=} F^N(M_\pi^N(k), \pi_k(M_\pi^N(k)))$ is the drift at time k if the empirical measure is $M_\pi^N(k)$. The first part is bounded with high probability using a Martingale argument (Lemma 4) and the second part is bounded using an integral formula.

Recall that $\bar{M}_\pi^N(t) \stackrel{\text{def}}{=} M_\pi^N\left(\left\lfloor \frac{t}{I(N)} \right\rfloor\right)$, i.e. $\bar{M}_\pi^N(kI(N)) = M_\pi^N(k)$ for $k \in \mathbb{N}$ and \bar{M}_π^N is piecewise constant and right-continuous. Let $\Delta_\pi^N(k)$ be the number of objects that change state between time slots k and $k+1$. Thus,

$$\|M_\pi^N(k+1) - M_\pi^N(k)\| \leq N^{-1} \sqrt{2} \Delta_\pi^N(k) \quad (30)$$

and thus

$$\|\hat{M}_\pi^N(t) - \bar{M}_\pi^N(t)\| \leq N^{-1} \sqrt{2} \Delta_\pi^N(k) \quad (31)$$

as well, with $k = \left\lfloor \frac{t}{I(N)} \right\rfloor$. Define

$$Z_\pi^N(k) = M_\pi^N(k) - M^N(0) - \sum_{j=0}^{k-1} F^N(M_\pi^N(j), \pi_j(M_\pi^N(j))) \quad (32)$$

and let $\hat{Z}_\pi^N(t)$ be the continuous, piecewise linear interpolation such that $\hat{Z}_\pi^N(kI(N)) = Z_\pi^N(k)$ for $k \in \mathbb{N}$. Recall that $A_\pi^N(t) \stackrel{\text{def}}{=} \pi_{\lfloor t/I(N) \rfloor}(M^N(\lfloor t/I(N) \rfloor)) - A_\pi^N(t)$ is the action taken by the controller at time $t/I(N)$. It follows from these definitions that:

$$\begin{aligned} \hat{M}_\pi^N(t) &= M_\pi^N(0) + \int_0^t \frac{1}{I(N)} F^N(\bar{M}_\pi^N(s), A_\pi^N(s)) ds + \hat{Z}_\pi^N(t) \\ &= M_\pi^N(0) + \int_0^t \frac{1}{I(N)} F^N(\hat{M}_\pi^N(s), A_\pi^N(s)) ds + \hat{Z}_\pi^N(t) \\ &\quad + \int_0^t \frac{1}{I(N)} \left[F^N(\bar{M}_\pi^N(s), A_\pi^N(s)) - F^N(\hat{M}_\pi^N(s), A_\pi^N(s)) \right] ds \end{aligned}$$

Using the definition of the semi-flow $\phi_t(m_0, A_\pi^N) = m_0 + \int_0^t f(\phi_s(m_0, A_\pi^N), A_\pi^N(s)) ds$, we get:

$$\begin{aligned} \hat{M}_\pi^N(t) - \phi_t(m_0, A_\pi^N) &= M_\pi^N(0) - m_0 + \hat{Z}_\pi^N(t) \\ &\quad + \int_0^t \frac{1}{I(N)} \left[F^N(\hat{M}_\pi^N(s), A_\pi^N(s)) - F^N(\phi_s(m_0, A_\pi^N), A_\pi^N(s)) \right] ds \\ &\quad + \int_0^t \left[\frac{1}{I(N)} F^N(\phi_s(m_0, A_\pi^N), A_\pi^N(s)) - f(\phi_s(m_0, A_\pi^N), A_\pi^N(s)) \right] ds \\ &\quad + \int_0^t \frac{1}{I(N)} \left[F^N(\bar{M}_\pi^N(s), A_\pi^N(s)) - F^N(\hat{M}_\pi^N(s), A_\pi^N(s)) \right] ds \end{aligned}$$

Applying Assumption (A2) to the third line, (A3) to the second and fourth lines, and Equation (31) to the fourth line leads to:

$$\begin{aligned} \left\| \hat{M}_\pi^N(t) - \phi_t(m_0, A_\pi^N) \right\| &\leq \left\| M_\pi^N(0) - m_0 \right\| + \left\| \hat{Z}_\pi^N(t) \right\| + L_1 \int_0^t \left\| \hat{M}_\pi^N(s) - \phi_s(m_0, A_\pi^N) \right\| ds \\ &\quad + I_0(N)t + \frac{\sqrt{2}L_1 I(N)}{N} \sum_{k=0}^{\lfloor \frac{t}{I(N)} \rfloor} \Delta_\pi^N(k) \end{aligned}$$

For all $N, \pi, T, b_1 > 0$ and $b_2 > 0$, define

$$\Omega_1 = \left\{ \omega \in \Omega : \sup_{0 \leq k \leq \frac{T}{I(N)}} \sum_{j=0}^k \Delta_\pi^N(j) > b_1 \right\}, \quad \Omega_2 = \left\{ \omega \in \Omega : \sup_{0 \leq k \leq \frac{T}{I(N)}} \left\| Z_\pi^N(k) \right\| > b_2 \right\} \quad (33)$$

Assumption (A1) implies conditions on the first and second order moment of $\Delta_\pi^N(k)$. Therefore by Lemma 3, this shows that for any $b_1 > 0$:

$$\mathbb{P}(\Omega_1) \leq \frac{TN^2}{b_1^2} \left[I_2(N) + \frac{I_1(N)^2}{I(N)^2} (T + I(N)) \right] \quad (34)$$

Moreover, we show in Lemma 4 that:

$$\mathbb{P}(\Omega_2) \leq 2S^2 \frac{T}{b_2^2} \left[2I_2(N) + I(N) [(I_0(N) + L_2)]^2 \right] \quad (35)$$

Now fix some $\epsilon > 0$ and let $b_1 = \frac{N\epsilon}{2\sqrt{2}L_1 I(N)}$, $b_2 = \epsilon/2$. For $\omega \in \Omega \setminus (\Omega_1 \cup \Omega_2)$ and for $0 \leq t \leq T$:

$$\begin{aligned} \left\| \hat{M}_\pi^N(t) - \phi_t(m_0, A_\pi^N) \right\| &\leq \left\| M_\pi^N(0) - m_0 \right\| + \epsilon + I_0(N)T \\ &\quad + L_1 \int_0^t \left\| \hat{M}_\pi^N(s) - \phi_s(m_0, A_\pi^N) \right\| ds \end{aligned}$$

By Grönwall's lemma:

$$\left\| \hat{M}_\pi^N(t) - \phi_t(m_0, A_\pi^N) \right\| \leq \left[\left\| M_\pi^N(0) - m_0 \right\| + \epsilon + I_0(N)T \right] e^{L_1 t} \quad (36)$$

and this is true for all $\omega \in \Omega \setminus (\Omega_1 \cup \Omega_2)$. We apply the union bound $\mathbb{P}(\Omega_1 \cup \Omega_2) \leq \mathbb{P}(\Omega_1) + \mathbb{P}(\Omega_2)$ which, with Eq.(34) and Eq.(35), concludes the proof.

The proof of Theorem 5 uses the following lemmas.

Lemma 3. *Let $(W_k)_{k \in \mathbb{N}}$ be a sequence of square integrable, non-negative random variables, adapted to a filtration $(\mathcal{F}_k)_{k \in \mathbb{N}}$, such that $W_0 = 0$ a.s. and for all $k \in \mathbb{N}$: $\mathbb{E}(W_{k+1} | \mathcal{F}_k) \leq \alpha$ and $\mathbb{E}(W_{k+1}^2 | \mathcal{F}_k) \leq \beta$. Then for all $n \in \mathbb{N}$ and $b > 0$:*

$$\mathbb{P} \left(\sup_{0 \leq k \leq n} (W_0 + \dots + W_k) > b \right) \leq \frac{n\beta + n(n+1)\alpha^2}{b^2} \quad (37)$$

Proof. Let $Y_n = \sum_{k=0}^n W_k$. It follows that $\mathbb{E}(Y_n) \leq \alpha n$ and

$$\mathbb{E}(Y_{n+1}^2) \leq \beta + 2n\alpha^2 + \mathbb{E}(Y_n^2)$$

from where we derive that

$$\mathbb{E}(Y_n^2) \leq n\beta + n(n+1)\alpha^2 \quad (38)$$

Now, because $W_{n+1} \geq 0$:

$$\mathbb{E}(Y_{n+1}^2 | \mathcal{F}_n) \geq (\mathbb{E}(Y_{n+1} | \mathcal{F}_n))^2 = (Y_n + \mathbb{E}(W_{n+1} | \mathcal{F}_n))^2 \geq Y_n^2$$

thus Y_n^2 is a non-negative sub-martingale and by Kolmogorov's inequality:

$$P\left(\sup_{0 \leq k \leq n} Y_k > b\right) = P\left(\sup_{0 \leq k \leq n} Y_k^2 > b^2\right) \leq \frac{\mathbb{E}(Y_n^2)}{b^2}$$

Together with Eq.(38) this concludes the proof. \square

Lemma 4. Define Z_π^N as in Eq.(32). For all $N \geq 2$, $b > 0$, $T > 0$ and all policy π :

$$\mathbb{P}\left(\sup_{0 \leq k \leq \lfloor \frac{T}{T(N)} \rfloor} \|Z_\pi^N(k)\| > b\right) \leq 2S^2 \frac{T}{b^2} \left[2I_2(N) + I(N) [(I_0(N) + L_2)]^2\right]$$

Proof. The proof is inspired by the methods in [1]. For fixed N and $h \in \mathbb{R}^S$, let

$$L_k = \langle h, Z_\pi^N(k) \rangle$$

By the definition of Z^N , L_k is a martingale w.r. to the filtration $(\mathcal{F}_k)_{k \in \mathbb{N}}$ generated by M_π^N . Thus

$$\mathbb{E}\left((L_{k+1} - L_k)^2 \middle| \mathcal{F}_k\right) = \mathbb{E}\left(\langle h, M_\pi^N(k+1) - M_\pi^N(k) \rangle^2 \middle| \mathcal{F}_k\right) + \langle h, F^N(M_\pi^N(k), \pi_k(M_\pi^N(k))) \rangle^2$$

By Assumption (A2):

$$|\langle h, F^N(M_\pi^N(k), \pi(M_\pi^N(k))) \rangle| \leq (I_0(N) + L_2) I(N) \|h\|$$

Thus, using Eq.(30) and Assumption (A1):

$$\begin{aligned} \mathbb{E}\left((L_{k+1} - L_k)^2 \middle| \mathcal{F}_k\right) &\leq \|h\|^2 \left[N^{-2} 2\mathbb{E}(\Delta_\pi^N(k)^2 \middle| \mathcal{F}_k) + [(I_0(N) + L_2) I(N)]^2\right] \\ &\leq \|h\|^2 \left[2I(N)I_2(N) + [(I_0(N) + L_2) I(N)]^2\right] \end{aligned}$$

We now apply Kolmogorov's inequality for martingales and obtain

$$\mathbb{P}\left(\sup_{0 \leq k \leq n} L_k > b\right) \leq \frac{n}{b^2} \|h\|^2 \left[2I(N)I_2(N) + [(I_0(N) + L_2) I(N)]^2\right]$$

Let Ξ_h be the set of $\omega \in \Omega$ such that $\sup_{0 \leq k \leq n} \langle h, Z_\pi^N(k) \rangle \leq b$ and let $\Xi := \bigcap_{h=\pm \vec{e}_i, i=1 \dots S} \Xi_h$ where \vec{e}_i is the i th vector of the canonical basis of \mathbb{R}^S . It follows that, for all $\omega \in \Xi$ and $0 \leq k \leq n$ and $i = 1 \dots S$: $|\langle Z_\pi^N(k), \vec{e}_i \rangle| \leq b$. This means that for all $\omega \in \Xi$: $\|Z_\pi^N(k)\| \leq \sqrt{S}b$. By the union bound applied to the complement of Ξ , we have

$$1 - \mathbb{P}(\Xi) \leq 2S \frac{n}{b^2} \left[I(N)I_2(N) + [(I_0(N) + L_2) I(N)]^2\right]$$

Thus we have shown that, for all $b > 0$:

$$\mathbb{P}\left(\sup_{0 \leq k \leq n} \|Z_\pi^N(k)\| > \sqrt{S}b\right) \leq 2S \frac{nI(N)}{b^2} \left[I_2(N) + I(N) [(I_0(N) + L_2)]^2\right]$$

which, by changing b to b/\sqrt{S} , shows the result. \square

5.4 Proof of Theorem 6

We use the same notation as in the proof of Theorem 5. By definition of V^N , v and the time horizons:

$$\begin{aligned} V_\pi^N(M^N(0)) - \mathbb{E}(v_{A_\pi^N}(m_0)) &= \mathbb{E} \left(\int_0^{H^N I(N)} r(\bar{M}_\pi^N(s), A_\pi^N(s)) - r(m_{A_\pi^N}(s), A_\pi^N(s)) ds \right) \\ &\quad - \mathbb{E} \left(\int_0^T r(m_{A_\pi^N}(s), A_\pi^N(s)) ds \right) \end{aligned}$$

The latter term is bounded by $I(N) \|r\|_\infty$. Let $\epsilon > 0$ and $\Omega_0 = \Omega_1 \cup \Omega_2$ where Ω_1, Ω_2 are as in the proof of Theorem 5. Thus $\mathbb{P}(\Omega_0) \leq \frac{J(N,T)}{\epsilon^2}$ and, using the Lipschitz continuity of r in m (with constant K_r):

$$\begin{aligned} |V_\pi^N(M^N(0)) - \mathbb{E}[v_{A_\pi^N}(m_0)]| &\leq I(N) \|r\|_\infty + \frac{2 \|r\|_\infty J(N,T)}{\epsilon^2} + \\ &\quad K_r \mathbb{E} \left[1_{\omega \notin \Omega_0} \int_0^T \|\bar{M}_\pi^N(s) - m_{A_\pi^N}(s)\| ds \right] \end{aligned}$$

For $\omega \notin \Omega_0$ and $s \in [0, T]$: $\int_0^T \|\bar{M}_\pi^N(s) - \hat{M}_\pi^N(s)\| ds \leq \frac{\epsilon I(N)}{2L_1}$ and, by Eq.(36), $\int_0^T \|\hat{M}_\pi^N(s) - m_{A_\pi^N}(s)\| ds \leq (\|M^N(0) - m_0\| + I_0(N)T + \epsilon) \frac{e^{L_1 T} - 1}{L_1}$ thus

$$|V_\pi^N(M^N(0)) - \mathbb{E}[v_{A_\pi^N}(m_0)]| \leq B_\epsilon(N, \|M^N(0) - m_0\|) \quad (39)$$

where

$$B_\epsilon(N, \delta) \stackrel{\text{def}}{=} I(N) \|r\|_\infty + K_r (\delta + I_0(N)T + \epsilon) \frac{e^{L_1 T} - 1}{L_1} + \frac{K_r I(N)}{2L_1} \epsilon + \frac{2 \|r\|_\infty J(N,T)}{\epsilon^2}$$

This holds for every $\epsilon > 0$, thus

$$|V_\pi^N(M^N(0)) - \mathbb{E}[v_{A_\pi^N}(m_0)]| \leq B(N, \|M^N(0) - m_0\|) \quad (40)$$

where $B(N, \delta) \stackrel{\text{def}}{=} \inf_{\epsilon > 0} B_\epsilon(N, \delta)$. By direct calculus, one finds that $\inf_{\epsilon > 0} (a\epsilon + b/\epsilon^2) = 3/2^{2/3} a^{2/3} b^{1/3}$ for $a > 0, b > 0$, which gives the required formula for $B(N, \delta)$.

5.5 Proof of Theorem 3

Let $\bar{\alpha}^N$ be the right-continuous function constant on the intervals $[kI(N); (k+1)I(N))$ such that $\bar{\alpha}^N(s) = \alpha(s)$. $\bar{\alpha}^N$ can be viewed as a policy independent of m . Therefore, by Theorem 5, on the set $\Omega \setminus (\Omega_1 \cup \Omega_2)$, for every $t \in [0; T]$:

$$\|\hat{M}_\alpha(t) - \phi_t(m_0, \alpha)\| \leq [\|M^N(0) - m_0\| + I_0(N)T + \epsilon] e^{L_1 T} + u(t)$$

with $u(t) \stackrel{\text{def}}{=} |\phi_t(m_0, \bar{\alpha}^N) - \phi_t(m_0, \alpha)|$. We have

$$\begin{aligned} u(t) &\leq \int_0^t |f(\phi_s(m_0, \alpha), \alpha(s)) - f(\phi_s(m_0, \bar{\alpha}^N), \bar{\alpha}^N(s))| ds \\ &\leq \int_0^t K (\|\phi_s(m_0, \alpha) - \phi_s(m_0, \bar{\alpha}^N)\| + d(\alpha(s), \bar{\alpha}^N(s))) ds \\ &\leq K \int_0^t u(s) ds + K d_1 \end{aligned}$$

where $d_1 \stackrel{\text{def}}{=} \int_0^T \|\alpha(t) - \bar{\alpha}^N(t)\| dt$. Therefore, using Grönwall's inequality, we have $u(t) \leq K d_1 e^{K T}$. By Lemma 5, this shows Eq.(13). The rest of the proof is as for Theorem 6.

Lemma 5. *If α is a piecewise Lipschitz continuous action function on $[0; T]$, of constant K_α , and with at most p discontinuity points, then*

$$\int_0^T d(\alpha(t), \bar{\alpha}^N(t)) dt \leq TI(N) \left(\frac{K_\alpha}{2} + 2(1 + \min(1/I(N), p)) \|\alpha\|_\infty \right).$$

Proof of lemma 5. Let first assume that $T = kI(N)$. The left handside $d_1 = \int_0^T d(\alpha(t), \bar{\alpha}^N(t)) dt$ can be decomposed on all intervals $[iI(N), (i+1)I(N))$:

$$d_1 = \sum_{i=0}^{\lfloor T/I(N) \rfloor} \int_{iI(N)}^{(i+1)I(N)} \|\alpha(s) - \bar{\alpha}^N(s)\| ds \leq \sum_{i=0}^{\lfloor T/I(N) \rfloor} \int_{iI(N)}^{(i+1)I(N)} \|\alpha(s) - \alpha(iI(N))\| ds$$

If α has no discontinuity point on $[iI(N), (i+1)I(N))$, then

$$\int_{iI(N)}^{(i+1)I(N)} d(\alpha(s), \alpha(iI(N))) ds \leq \int_0^{I(N)} K_\alpha s ds \leq K_\alpha 2I(N)^2$$

If α has one or more discontinuity points on $[iI(N), (i+1)I(N))$, then

$$\int_{iI(N)}^{(i+1)I(N)} d(\alpha(s), \alpha(iI(N))) ds \leq \int_{iI(N)}^{(i+1)I(N)} 2\|\alpha\|_\infty ds \leq 2\|\alpha\|_\infty I(N)$$

There are at most $\min(1/I(N), p)$ intervals $[iI(N), (i+1)I(N)]$ that have discontinuity points which shows that

$$d_1 \leq TI(N) \left(\frac{K_\alpha}{2} + \min(1/I(N), p) 2\|\alpha\|_\infty \right).$$

If $T \neq kI(N)$, then $T = kI(N) + t$ with $0 < t < I(N)$. Therefore, there is an additional term of $\int_{kI(N)}^{kI(N)+t} d(\alpha(s), \bar{\alpha}^N(s)) ds \leq 2\|\alpha\|_\infty I(N)$. \square

5.6 Proof of Theorem 2

This theorem is a direct consequence of Theorem 3 and Theorem 6. We do the proof for almost sure convergence, the proof for convergence in probability is similar. To prove the theorem we prove

$$\limsup_{N \rightarrow \infty} V_*^N(M^N(0)) \leq v_*(m_0) \leq \liminf_{N \rightarrow \infty} V_*^N(M^N(0)) \quad (41)$$

- Let $\epsilon > 0$ and $\alpha(\cdot)$ be an action function such that $v_\alpha(m_0) \geq v_*(m_0) - \epsilon$ (such an action is called ϵ -optimal). Theorem 3 shows that $\lim_{N \rightarrow \infty} V_\alpha^N(M^N(0)) = v_\alpha(m_0) \geq v_*(m_0) - \epsilon$ a.s. This shows that $\liminf_{N \rightarrow \infty} V_*^N(M^N(0)) \geq \lim_{N \rightarrow \infty} V_\alpha^N(M^N(0)) \geq v_*(m_0) - \epsilon$; this holds for every $\epsilon > 0$ thus $\liminf_{N \rightarrow \infty} V_*^N(M^N(0)) \geq v_*(m_0)$ a.s., which establishes the second inequality in Eq.(41), on a set of probability 1.
- Let $B(N, \delta)$ be as in Theorem 6, $\epsilon > 0$ and π^N such that $V_*^N(M^N(0)) \leq V_{\pi^N}^N(M^N(0)) + \epsilon$. By Theorem 6, $V_{\pi^N}^N(M^N(0)) \leq \mathbb{E} \left(v_{A_{\pi^N}^N}(m_0) \right) + B(N, \delta^N) \leq v_*(m_0) + B(N, \delta^N)$ where $\delta^N \stackrel{\text{def}}{=} \|M^N(0) - m_0\|$. Thus $V_*^N(M^N(0)) \leq v_*(m_0) + B(N, \delta^N) + \epsilon$. If further $\delta^N \rightarrow 0$ a.s. it follows that $\limsup_{N \rightarrow \infty} V_*^N(M^N(0)) \leq v_*(m_0) + \epsilon$ a.s. for every $\epsilon > 0$, thus $\limsup_{N \rightarrow \infty} V_*^N(M^N(0)) \leq v_*(m_0)$ a.s.

6 Conclusion and Perspectives

There are several natural questions arising from this work. One concerns the convergence of optimal policies. Optimal policies π_*^N of a stochastic systems with N objects may not be unique, they may also exhibit thresholds and therefore be discontinuous. This implies that $M_{\pi_*^N}^N$ and $V_{\pi_*^N}^N$ will not

converge in general. In some particular cases, such as the best response dynamics studied in [10], limit theorems can nevertheless be obtained, at the cost of a much greater complexity. In full generality however, this problem is still open and definitely deserves further investigations.

The second question concerns the time horizon. In this paper we have focused on the finite horizon case. Actually, most results and in particular theorems 2 and 3, remain valid with an infinite horizon with discount. The main argument that makes everything work in the discounted case is the following. When the rewards $r(s, a)$ are bounded, for a given discount $\beta < 1$ and a given $\varepsilon > 0$, it is possible to find a finite time horizon T such that the expected discounted value of a policy π can be decomposed into the value over time T plus a term less than ε :

$$\mathbb{E} \sum_{t \geq 0} \beta^t r(M^N(t), \pi(M^N(t))) \leq \mathbb{E} \sum_{t=0}^T \beta^t r(M^N(t), \pi(M^N(t))) + \varepsilon.$$

Therefore, the main result of this paper, which states that a policy π that is optimal in the mean field limit is near-optimal for the finite system with N objects, also holds in the infinite horizon discounted case.

As for the infinite horizon without discount or average reward cases, convergence of the value when N goes to infinity is not guaranteed in general. Finding natural assumptions under which convergence holds is also one of our goals for the future.

References

- [1] M. Benaïm. Dynamics of stochastic approximation algorithms. *Séminaire de Probabilités XXXIII. Lecture Notes in Math*, 1709:1–68, 1999.
- [2] M. Benaïm and J.-Y. Le Boudec. A Class Of Mean Field Interaction Models for Computer and Communication Systems. *Performance Evaluation*, 65(11-12):823–838, 2008.
- [3] M. Benaïm and J. Weibull. Deterministic approximation of stochastic evolution in games: a generalization. Technical report, mimeo, 2003.
- [4] A. Benveniste, P. Priouret, and M. Métivier. *Adaptive algorithms and stochastic approximations*. Springer-Verlag, 1990.
- [5] V. Berten and B. Gaujal. Grid brokering for batch allocation using indexes. In Springer, editor, *Network Control and Optimization*, volume 4465 of *LNCS*, 2007.
- [6] W. Chen, D. Huang, A.A. Kulkarni, J. Unnikrishnan, Q. Zhu, P. Mehta, S. Meyn, and A. Wierman. Approximate dynamic programming using fluid and diffusion approximations with applications to power management. In *Decision and Control, 2009 held jointly with the 2009 28th Chinese Control Conference. CDC/CCC 2009. Proceedings of the 48th IEEE Conference on*, pages 3575–3580. IEEE, 2009.
- [7] R. Cole. Initial studies on worm propagation in manets for future army combat systems. Technical report, Pentagon Reports, 2004.
- [8] DP De Farias and B. Van Roy. The linear programming approach to approximate dynamic programming. *Operations Research*, 51(6):850–865, 2003.
- [9] N. Gast and B. Gaujal. A mean field approach for optimization in discrete time. *Discrete Event Dynamic Systems*, 21:63–101, 2011.
- [10] Z. Gorodeisky. Deterministic approximation of best-response dynamics for the Matching Pennies game. *Games and Economic Behavior*, 66(1):191–201, 2009.
- [11] M. Huang, P. E. Caines, and R. P. Malhame. Large population stochastic dynamic games: closed-loop mckean-vlasov systems and the nash certainty equivalence principle. *Commun. Inform. Systems*, 6(3):221–252, 2006.

- [12] M. Huang, P. E. Caines, and R. P. Malhame. Nash certainty equivalence in large population stochastic dynamic games: connections with the physics of interacting particle systems. In *45th IEEE Conference on Decision and Control*, pages 4921–4926, San Diego, 2006.
- [13] M. Huang, P. E. Caines, and R. P. Malhame. Large-population cost-coupled lqg problems with nonuniform agents: individual-mass behavior and decentralized e-nash equilibria. *IEEE Transactions Automatic Control*, 52, 2007.
- [14] M. Huang, P. E. Caines, and R. P. Malhame. Social optima in mean field lqg control: centralized and decentralized strategies. In *47th Allerton Conference*, 2009.
- [15] M.H.R. Khouzani, S. Sarkar, and E. Altman. Maximum damage malware attack in mobile wireless networks. In *IEEE Infocom*, San Diego, 2010.
- [16] D. Kondo, B. Javadi, P. Malecot, F. Cappello, and D. Anderson. Cost-benefit analysis of cloud computing versus desktop grids. In *18th International Heterogeneity in Computing Workshop*, Rome, 2009.
- [17] T. Kurtz. Solutions of ordinary differential equations as limits of pure jump markov processes. *Journal of Applied Probability*, pages 49–58, 1970.
- [18] J.-M. Lasry and P.-L. Lions. Mean field games. *Japan. Journal Math.*, 2007.
- [19] J.Y. Le Boudec, D. McDonald, and J. Mundinger. A generic mean field convergence result for systems of interacting objects. In *Quantitative Evaluation of Systems, 2007. QEST 2007. Fourth International Conference on the*, pages 3–18, 2007.
- [20] C. H. Papadimitriou and J. N. Tsitsiklis. The complexity of optimal queuing network control. *Math. Oper. Res.*, 24(2):292–305, 1999.
- [21] M.L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. Series in Probability and Mathematical Statistics. Wiley, 1994.
- [22] S. Tanachaiwiwat and A. Helmy. Vaccine: War of the worms in wired and wireless networks. In *IEEE INFOCOM*, 2006.
- [23] H. Tembine, J.-Y. Le Boudec, R. El-Azouzi, and E. Altman. Mean field asymptotic of markov decision evolutionary games and teams. In *Gamenets*, 2009.
- [24] H. Tembine, P. Vilanova, and M. Debbah. Noisy mean field stochastic games with network applications. Technical Report 342-P-11867-83, Supélec, 2010.
- [25] J.N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690, 1997.



Centre de recherche INRIA Grenoble – Rhône-Alpes
655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Centre de recherche INRIA Bordeaux – Sud Ouest : Domaine Universitaire - 351, cours de la Libération - 33405 Talence Cedex
Centre de recherche INRIA Lille – Nord Europe : Parc Scientifique de la Haute Borne - 40, avenue Halley - 59650 Villeneuve d'Ascq
Centre de recherche INRIA Nancy – Grand Est : LORIA, Technopôle de Nancy-Brabois - Campus scientifique
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex
Centre de recherche INRIA Paris – Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex
Centre de recherche INRIA Rennes – Bretagne Atlantique : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex
Centre de recherche INRIA Saclay – Île-de-France : Parc Orsay Université - ZAC des Vignes : 4, rue Jacques Monod - 91893 Orsay Cedex
Centre de recherche INRIA Sophia Antipolis – Méditerranée : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex

Éditeur
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)
<http://www.inria.fr>
ISSN 0249-6399